



다음과 같이 데이터가 주어졌다고 생각해 보자. feature가 output에 선형적으로 관여한다고 생각하면 다음과 같이 표현할 수 있다.

$$y \approx \theta_0 + \theta_1 x_1$$

우리의 목표는 위의 값의 좌우가 최대한 일치하도록 θ_0, θ_1 을 결정해야 한다. 그러기 위해서는 다음의 cost-function의 값을 최소화시켜 주어야 한다.

$$\frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = J(\theta)$$

예를 들어 값이 다르면 다른수록 $(h_{\theta}(x^{(i)}) - y^{(i)})^2$ 이 커지기 때문이다.

이 $J(\theta)$ 값이 작아지도록 θ 를 설정하기 위해 우리는 gradient descent를 사용한다.

gradient descent란 $J(\theta)$ 에 대하여 $\frac{\partial J(\theta)}{\partial \theta_i}$ 를 구하고, 만약 0보다 이 값이 작을 경우 θ_i 를 늘렸을 때 값이 작아지므로 θ_i 를 + 방향으로 옮긴다. 또한 0보다 클 경우에는 θ_i 를 늘렸을 때 값이 커지므로 θ_i 를 - 방향으로 옮긴다. 그러므로 θ_i 를 다음과 같이 갱신해야 할 수 있다.

$$\theta_i \leftarrow \theta_i - \alpha \frac{\partial J(\theta)}{\partial \theta_i}$$

원래의 $J(\theta)$ 와 갱신된 θ 들을 넣은 $J(\theta)$ 가 거의 다르지 않을 경우 갱신을 종료한다.

2. multiple-variable의 경우 $h_0(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 \dots \theta_n x_n$ 이된다.

\therefore Cost function은

$$J(\theta) = J(\theta_0 \dots \theta_n) = \frac{1}{2m} \sum_{i=1}^m (\theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} \dots \theta_n x_n^{(i)} - y^{(i)})^2$$

이된다.

위와 같이 정의할수 있는 이유는 $x_0, x_1 \dots x_n$ 이 y 에 대해 선형적으로 관여한다고 생각했기 때문이다. 따라서

$$y \approx \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 \dots \theta_n x_n = h_0(x)$$

가 되고 좌변과 우변 차이가 가장 작은 $\theta_0 \dots \theta_n$ 을 구하면 되는 것이다.

그러기 위해서는 $|\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 \dots \theta_n x_n - y|$ 를 작게 만들어야 한다.

하지만 절댓값이 미분은 조금 복잡할수 있으므로 $(\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 \dots \theta_n x_n - y)^2$ 을 가장 최소가 되도록 설정해두면된다. \therefore 위의 $J(\theta)$ 는 cost function에 적절한 함수임을 알수 있다.

3. $h_\theta(x)$ 과 y 가 비슷해지기 위해서

cost function인 $\frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 = J(\theta)$ 가 최소가 되어야 한다.

$J(\theta)$ 값이 작아지도록 θ 를 설정하기 위해 우리는 gradient descent를 사용한다.

gradient descent란 $J(\theta)$ 에 대하여 $\frac{\partial J(\theta)}{\partial \theta_i}$ 를 구하고, 만약 0보다 이 값이 작을

경우 θ_i 를 늘렸을 때 값이 작아지므로 θ_i 를 +방향으로 움직여 준다. 또한 0보다 클 경우에는

θ_i 를 늘렸을 때 값이 작아지므로 θ_i 를 +방향으로 움직여 줄을 알 수 있다. 그러므로

θ 를 다음과 같이 갱신해야 적절하다고 생각할 수 있다.

$$\theta_i \leftarrow \theta_i - \alpha \frac{\partial J(\theta)}{\partial \theta_i}$$

α 값이 매우 작다면 근속이 극히 미미하여
외려다들 θ 값을 갱신할 것이다.

원래의 $J(\theta)$ 와 갱신된 θ 를 넣은 $J(\theta)$ 가 거의 다르지 않을 경우 갱신을 종료한다.